

Unterstützung datenintensiver Forschung am KIT – Aktivitäten, Dienste und Erfahrungen

Achim Streit

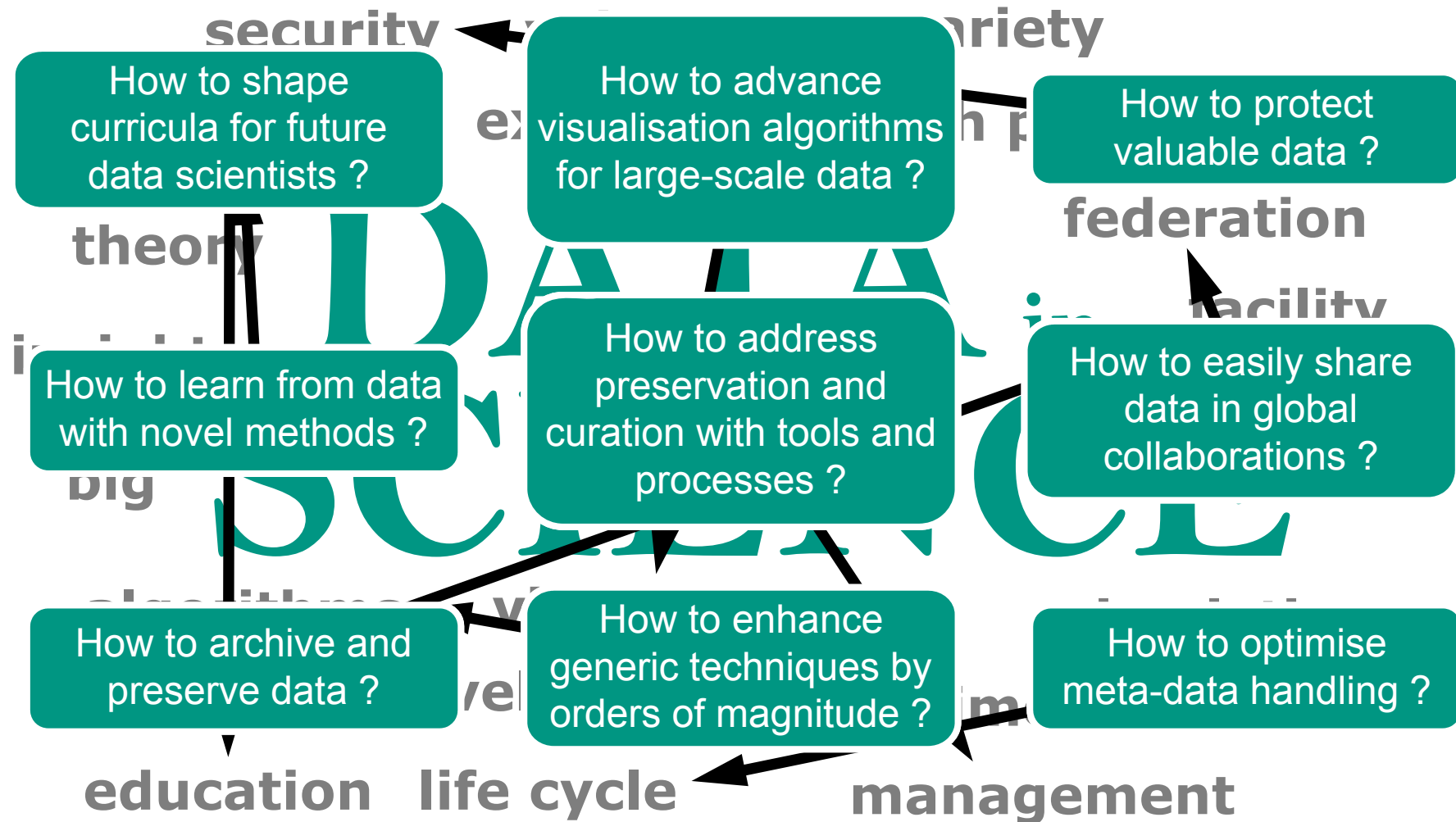
Steinbuch Centre for Computing (SCC)



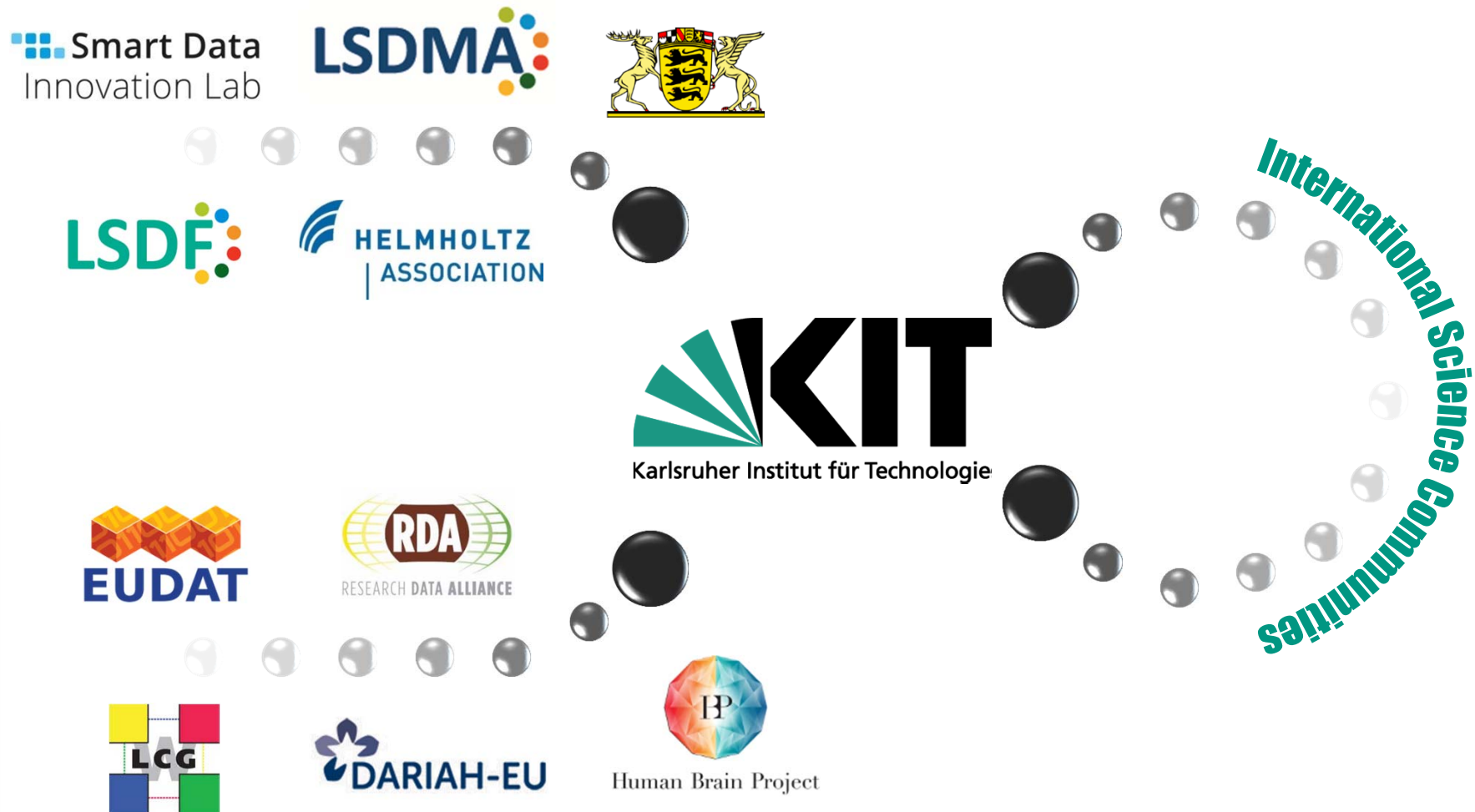
The Challenge of Big Data in Science

security volume variety
analysis exploration 4th paradigm
theory DATA in federation
insight SCIENCE facility
big
algorithms visualisation simulation
value velocity experiment veracity
education life cycle management

The Challenge of Big Data in Science



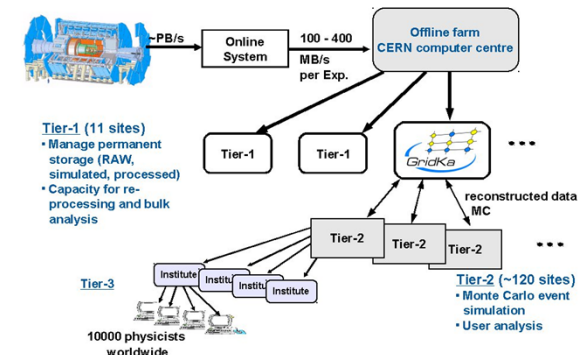
Collaboration is key



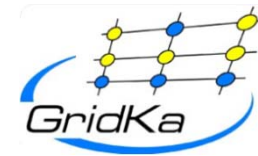
GridKa



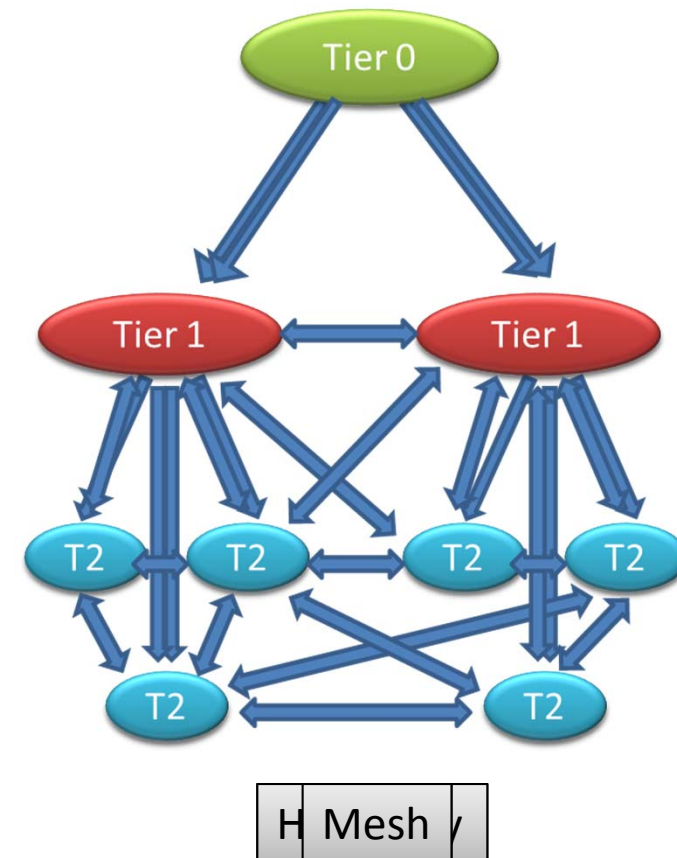
- German Tier-1 Center in WLCG
 - Supports all LHC experiments + Belle II + several small {HE/AP}P communities
 - Benchmarked reliability of 99.5%
- Resources
 - >10,000 cores, utilization >95%
 - Disk space: 12 PB, tape space: 17 PB
 - 6x10 Gbit/s network connectivity
- 14% of LHC data permanently stored at GridKa
- Serves > 20 T2 centres in 6 countries
- Services
 - File transfer
 - Regional workload management, file catalog
- Annual international GridKa School
- Global Grid User Support (GGUS) for WLCG



GridKa Experiences



- Evolving demands and usage patterns
 - No common workflows
- Hardware commodity, software not
- Hierarchical storage with tape challenging
- On-site experiment representation highly useful
- Adoption of grid computing outside of HEP rare
- Data access a central issue
 - Random data access by compute jobs
 - Reprocessing



Courtesy of Ian Bird, CERN

Large Scale Data Facility

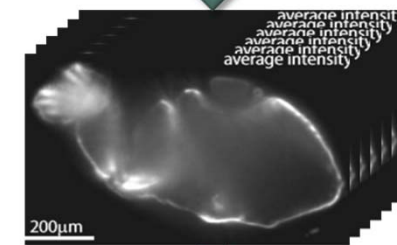
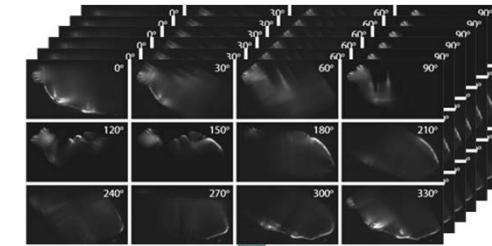


■ Main goals

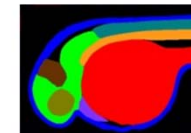
- Provision of storage for multiple research groups
 - Systemsbiology, climatology, synchrotron research/materials science, humanities, geo/earth science, ...
- Basis for BW-wide data services

■ Resources and access

- 6 PB of on-line storage (ext. to > 10 PB)
- 6 PB of archival storage
- 100 GbE connection between LSDF@KIT and U-Heidelberg
- Hadoop analysis cluster of 58*8 cores
- Connection to HPC clusters
- Jointly funded by Helmholtz Association and state of Baden-Württemberg



Model



LSDF Experiences



- High demand for storage, analysis and archival
- Research groups vary in
 - Research topics (from genetic sequencing to geophysics)
 - Size
 - IT expertise
 - Need for services and protocols
- Important needs common to many user groups
 - Sharing data with other groups
 - Data safety and preservation
 - 'Consulting'
- Many small groups depend on LSDF

bwLSDF – Services for users in the state of Baden-Württemberg



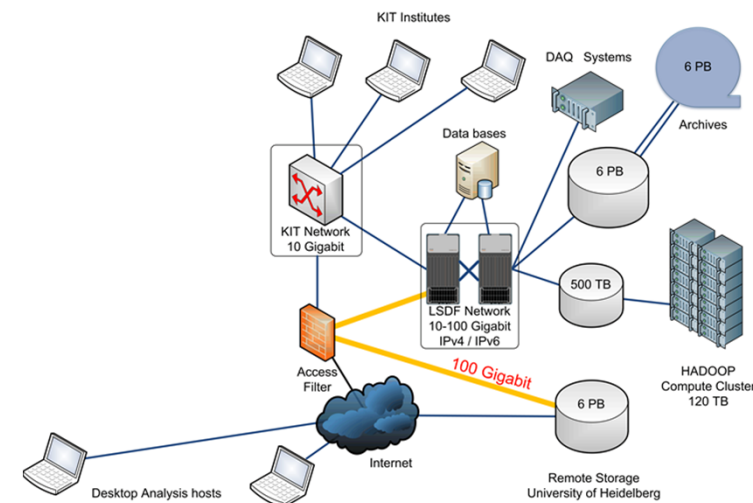
■ bwSync&Share

- For employees and students at Universities in BaWü
- Dropbox-like data storage
- KIT cloud – on-premise solution
- Access via webbrowser, clients for windows, linux and mobile phones
- Web based SAML authentication (Shibboleth) - bwIDM
- Data can be shared with users not registered for the service
- Service started on 1.1.2014
- Available to 450.000 users

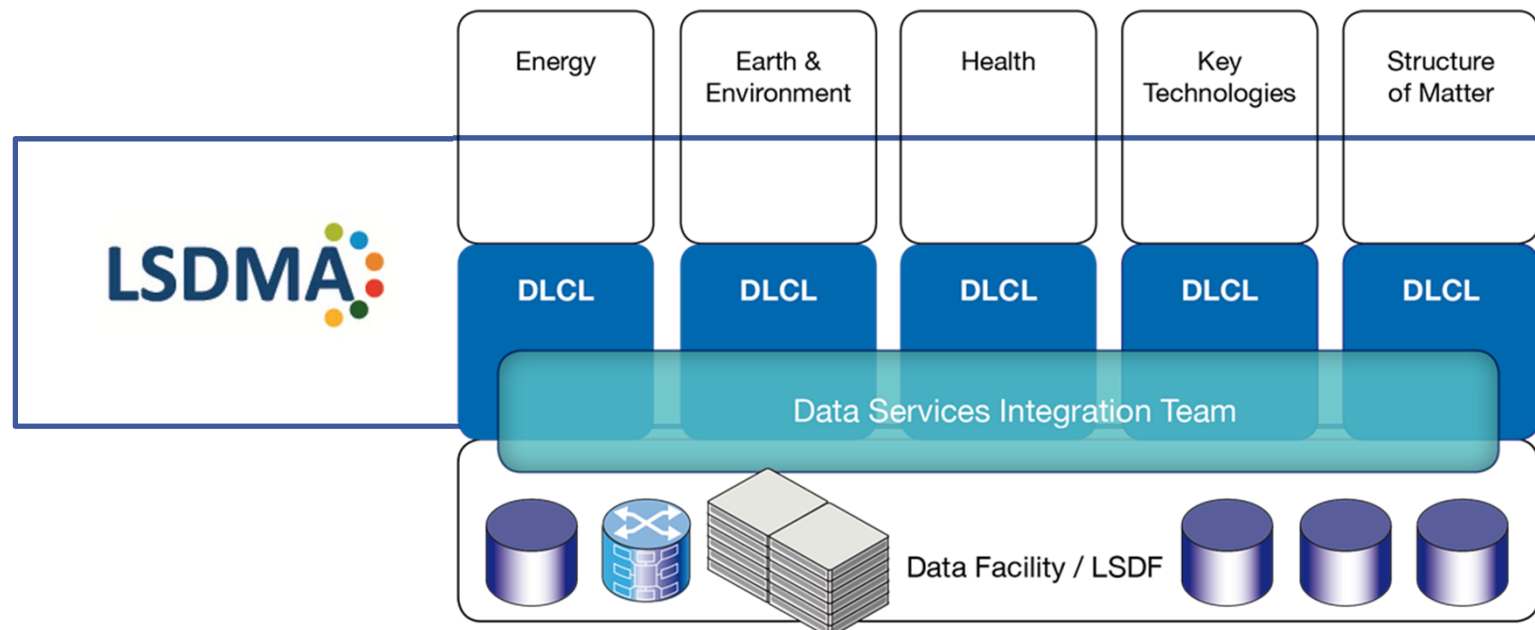


■ bwFileStorage

- Via SFTP/SCP
- SAML authentication



Large-Scale Data Management and Analysis – Dual Approach



Data Life Cycle Labs

Joint R&D with scientific user communities

- Optimization of the data life cycle
- Community-specific data analysis tools and services

Data Services Integration Team

Generic methods R&D

- Data analysis tools and services common to several DLCLs
- Interface between federated data infrastructures and DLCLs/communities

Facts and Figures



- Helmholtz portfolio extension
- Initial project duration: 2012-2016
- Partners:



ulm university universität



Universität Hamburg



- Sustainability
 - Inclusion of activities into Helmholtz programme-oriented funding (PoF) “Supercomputing & Big Data” from 2015 onwards
 - Cross-programme initiative with other Helmholtz research fields
- Annual international symposium

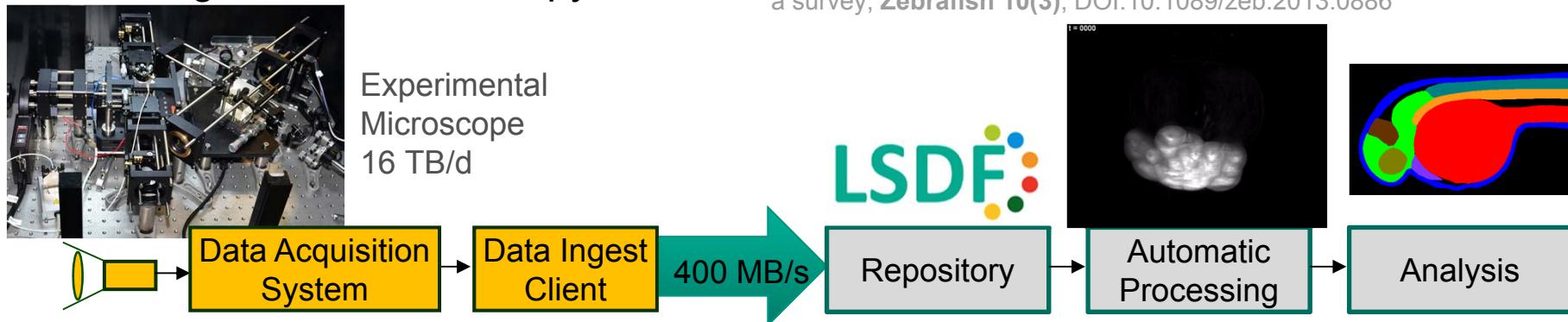


Selected Scientific Highlights

DLCL Key Technologies (KIT, U-Heidelberg, U-Dresden)

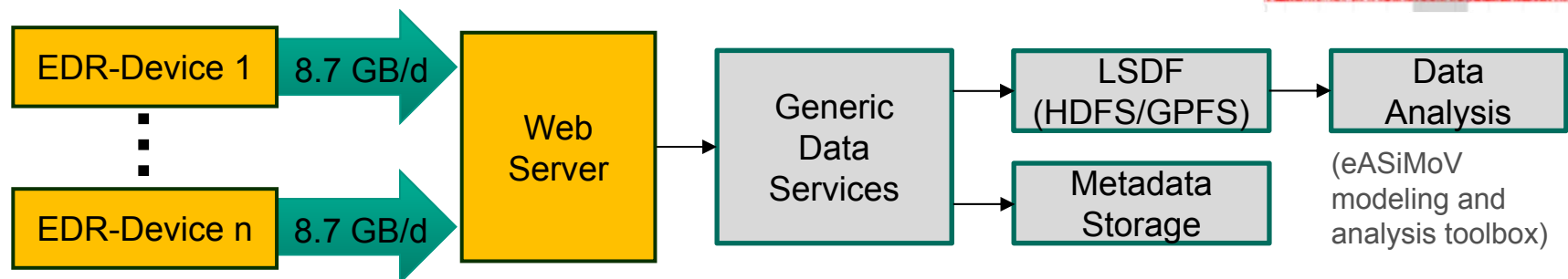
Light Sheet Microscopy

Mikut et al. (2013), Automated processing of zebrafish imaging data - a survey, *Zebrafish* 10(3), DOI:10.1089/zeb.2013.0886



DLCL Energy (KIT, U-Ulm)

- Secure data management for tech. & eco. data analysis
- Electrical Data Recorder: 3-phase voltage measurement



LSDMA Experiences

Communities differ in

- Previous knowledge
- Level of specification of the data life cycle
- Tools and services used

Needs driven by

- Increasing amount of data
- Cooperation between groups
- Policies
 - Open access/data
 - Long-term preservation

Within communities

- Focus on data analysis
- High fluctuation of computing experts

Lessons learned

- Interoperable AAI crucial
- Data privacy very challenging, both legally and technically
- Communities need evolution, not revolution
- Needs can be very specific

Re3data.org

© by KIT-BIB, Frank Scholze



- Goals
 - Linking existing research data repositories
 - In-depth analysis of quality requirements for research data repositories
 - Define a draft set of criteria for their quality assurance
- Currently 634 research data repositories from around the world covering all academic disciplines are listed
 - 586 of these are described using the re3data.org schema, <http://doi.org/10.2312/re3.005>



PANGAEA

Publishing Network for Geoscientific and Environmental Data



Subjects: Atmospheric Science and Oceanography Biology Geochemistry, Mineralogy and Crystallography
Geochemistry, Mineralogy and Crystallography Geology and Palaeontology Geology and Palaeontology Geophysics
Geophysics and Geodesy Geosciences (including Geography) Life Sciences Natural Sciences Oceanography

Content types: Archived data Audiovisual data Images Plain text Standard office documents

Countries: Germany

The information system PANGAEA is operated as an Open Access library aimed at archiving, publishing and distributing georeferenced data from earth system research. The system guarantees long-term availability of its content through a commitment of the operating institutions.

FUNDING



NETWORK



Smart Data innovation Lab: A Data Hub for Industry

 Smart Data
Innovation Lab



Governance

Data Innovation Communities

Industry 4.0



Energy



Smart Cities



Medicine



Working Group **Data Curation**



Working Group **Legal Affairs**

Cross Topic
Communities



Working Group
Facility Operation and Tools

